# Making California government sites more accessible to search engine users

## Implementing the Sitemap protocol

# Why we're talking today

- Sample search 1: ["richard abbott" california accountant license]

  (http://www.google.com/search?hl=en&lr=&safe=off&q=%22richard+abbott%22+california+accountant+license&btnG=Search)

  – What a search engine user does not find:

  http://www2.dca.ca.gov/pls/wllpub/WLLQRYNA$LCEV2.QueryView?P_LICENSE_NUMBER=2543&P_LTE_ID=781

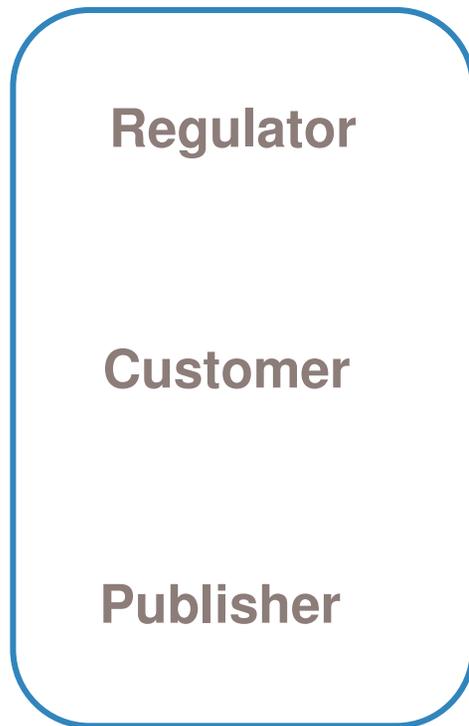  – Reason: Database uncrawlable

Google

# Why we're talking today

- Sample search 2: [california real estate appraisal licensing courses] (http://www.google.com/search?hl=en&lr=&safe=off&q=california+real+estate+appraisal+licensing+courses+&btnG=Search)

  - What a search engine user does not find: http://secure.dre.ca.gov/publicasp/CEStatutory.asp

  - Reason: Database uncrawlable

Google

# Why we're talking today

- Sample search 3: [cdec station search sensor type] (http://www.google.com/search?hl=en&q=cdec+station+search+sensor+type&btnG=Google+Search)

  – What a search engine user does not find: http://cdec.water.ca.gov/cgi-progs/staSearch; http://cdec.water.ca.gov/robots.txt

  – Reason: Site blocked by robots.txt (database itself also uncrawled)

Google

# Government's relationships with Google

## Government

**Google**

| Regulator | Policy | Corporate citizen |
|---|---|---|
| Customer | Enterprise | Vendor |
| Publisher | Content Partnerships | Distributor |

Google

# Agenda

Government information on the growing web

Sitemaps for search engines

Implementing the Sitemap protocol
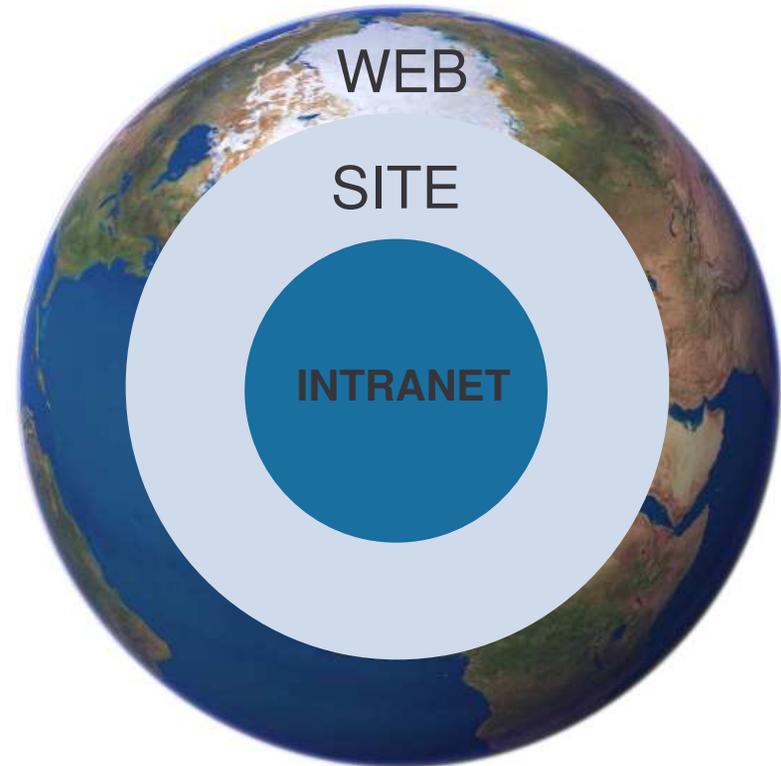
Google Webmaster Tools

Success stories

Q&A

Google

# Common concerns

- No direct cost

- Non-proprietary

- No security risk

- Public content only

Google

# What we're talking about

- **Intranet:** Not your intranet or internal information

- **Site:** Nor search within your public site

- **Web:** Making content on your site accessible to web search engine users



WEB

SITE

**INTRANET**
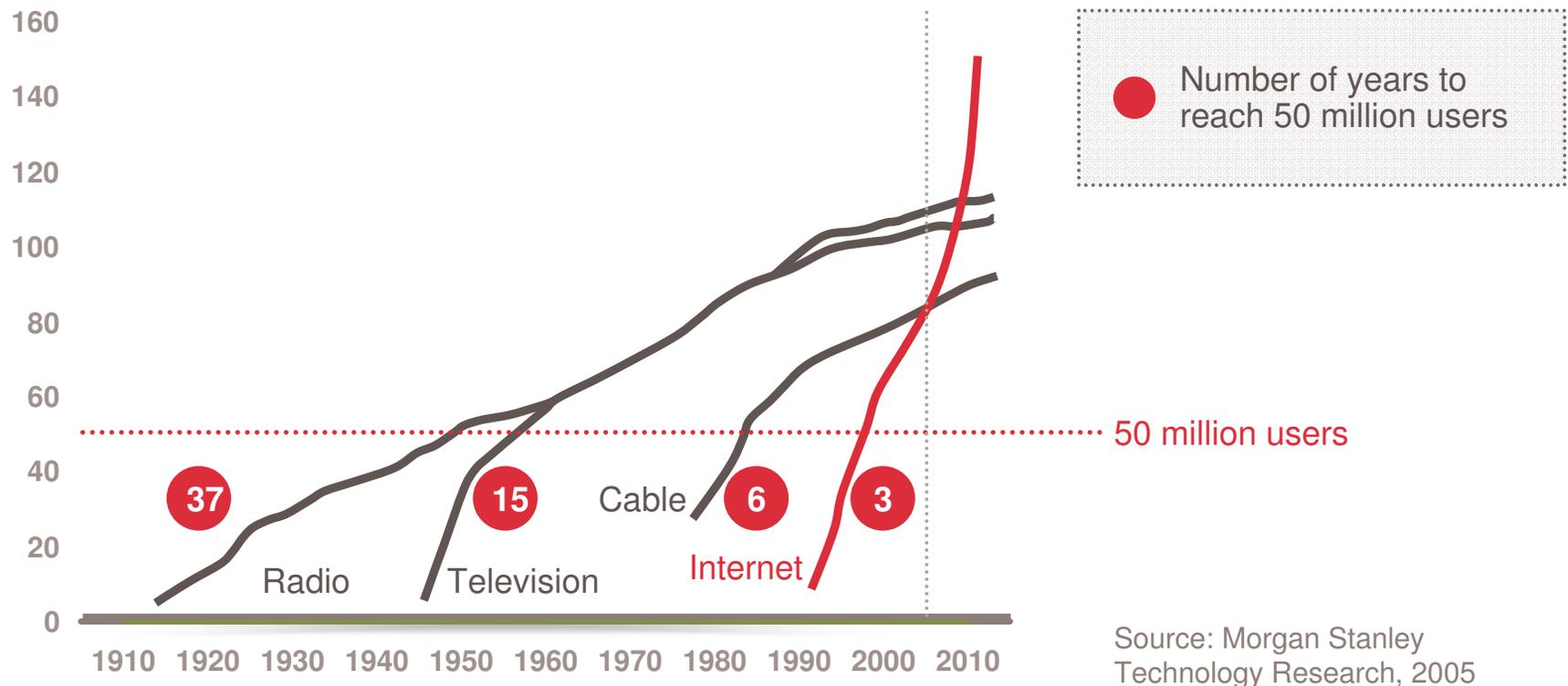
Google

# Web search vs. site search

## Supporting the two levels of search



| | Search scope | A segment of your public site's content |
|---|---|---|
| All of the open and accessible deep web | **Search scope** | A segment of your public site's content |
| Citizens and professionals | **User** | Professionals and citizens |
| Googlebot's crawling intervals | **Freshness** | Customizable |
| Limited by robots.txt, dynamic content. | **Crawling** | Limited by server capacity and cost |
| High-level stats | **Reporting tools** | More detailed, all facets |
| Free | **Cost** | Varies |

# The internet has come of age faster than previous media

Internet adoption by North American users/households



Number of years to reach 50 million users

50 million users

37 — Radio

15 — Television

Cable

6

3

Internet

1910 1920 1930 1940 1950 1960 1970 1980 1990 2000 2010

Source: Morgan Stanley Technology Research, 2005

Google

# US Internet user population is diversifying

73% of adult population is online

- Not just youngsters: 71% of baby boomers (50–64)

- Not just urban and suburban: 63% of rural residents

- Not just highly educated: 84% with "some college"



Source: Pew Internet & American Life Project, 2006

Google

# The growing web

## 7 million new pages every day

**17 LOC**

**.03 LOC**

**1 LOC**

Library at
Alexandria,
300 B.C.

Library of
Congress

**World Wide
Web**

Library of Congress (LOC)
**17M Books**

Source: Peter Lyman and Hal Varian, 2003

Google

# Not all information is created equal

The value of government content – a pillar of the web

**Government**

**Untrustworthy**                                                  **Trustworthy**

Google

# Citizens increasingly access government through search engines

## National Institutes of Health (nih.gov)

- 70% of unique users in July 2006 were referred by search engines (Google, Yahoo, MSN, AOL, Ask)



- Only 4% of unique users came directly to nih.gov

Source: ComScore, 2006

# And they expect to find everything

## The long tail of federal government information



Y-axis: Number of queries
X-axis: Obscurity of query

- irs refund
- fda drug approval
- veteran record deceased
- tree faller nps yosemite
- ftc hearing sherman act 1992

Google

# Search engines are the point of departure, government sites are the destination



**Federal**

Internal Revenue Service
DEPARTMENT OF THE TREASURY

NASA

**State**

virginia.gov

utah dot gov

**Localities**

King County

City of Dallas

Google Confidential

Google

# Government information on the growing web – recap

✓ Internet becoming dominant medium for accessing government

✓ Users value government information

✓ Users prefer to access government through search engines

Google Confidential

Google

# Growing deeper and more dynamic

Challenges to web crawling are growing and multiplying

- Outdated robots.txt crawling instructions

- Non-html links

- Content "hidden" behind search forms

- Server errors (crawler times out when fetching content)

- Orphaned URLs

- Rich media: audio, video

- Paid/premium content databases

**WEB**
Searchable

**DEEP WEB**
Not searchable

Google

# Crawlers cannot navigate search forms



**When crawled**



**Search results are invisible**

Google Confidential

Google

# The solution: Sitemaps

The Sitemap protocol enables a web publisher to proactively manage search engine crawling



"The launch of Sitemaps is significant because it allows for a single, easy way for websites to provide content and metadata to search engines"

—Tim Mayer, Senior Director of Product Management, Yahoo Search

"We are 100% behind this protocol - this kind of collaboration will help improve the search experience for all of our customers"

—Ken Moss, General Manager, Live Search

- Sitemap protocol developed by Google in June 2005 and released under Creative Commons License
- Adopted as an industry standard in November 2006: www.sitemaps.org

# Navigational sitemap

A browse index or sitemap enables a user to navigate throughout a site

Google Confidential

# Sitemaps for search engines

- HTML

- Simple text

- XML

Google

# Simple text sitemap

## A comprehensive list of URLs

http://www.firstgov.gov/index.shtml
http://www.firstgov.gov/About.shtml
http://www.firstgov.gov/Citizen/Services/Address_Changes.shtml
http://www.firstgov.gov/Topics/Parents_Adoptive.shtml
http://www.firstgov.gov/Government/State_Local/Ag_Environment.shtml
http://www.firstgov.gov/Citizen/Topics/Environment_Agriculture/Agriculture.shtml
http://www.firstgov.gov/Citizen/Facts/Facts_Agriculture.shtml
http://www.firstgov.gov/Agencies/Federal/Executive/Agriculture.shtml

Google

# XML sitemap

- A comprehensive list of URLs in XML

- Tagged with each URL's location, last modification, change frequency and priority

```xml
<?xml version="1.0" encoding="UTF-8"?>
<urlset xmlns="http://www.google.com/schemas/sitemap/0.84">

    <url>
        <loc>http://www.example.com/</loc>
        <lastmod>2005-01-01</lastmod>
        <changefreq>monthly</changefreq>
        <priority>0.8</priority>
    </url>
    <url>
        <loc>http://www.example.com/catalog?item=12&amp;desc=vacation_hawaii</loc>
        <changefreq>weekly</changefreq>
    </url>
    <url>
        <loc>http://www.example.com/catalog?item=73&amp;desc=vacation_new_zealand</loc>
        <lastmod>2004-12-23</lastmod>
        <changefreq>weekly</changefreq>
    </url>
    <url>
        <loc>http://www.example.com/catalog?item=74&amp;desc=vacation_newfoundland</loc>
        <lastmod>2004-12-23T18:00:15+00:00</lastmod>
        <priority>0.3</priority>
    </url>
    <url>
        <loc>http://www.example.com/catalog?item=83&amp;desc=vacation_usa</loc>
        <lastmod>2004-11-23</lastmod>
    </url>
</urlset>
```

Google